# NCBI RefSeq Resources for Plant Genomics

Anjana R Vatsan(raina@nih.gov)

Functional Genomics Workshop

PAG XXVIII January 13, 2020

U.S. National Library of Medicine
*National Center for Biotechnology Information*

# Also from NCBI!

| Day | Time | Topic |
| --- | --- | --- |
| Monday | 12:50 pm – 3:00 pm<br>*Pacific Salon 1* | **NCBI Genome Resources Workshop** |
| Tuesday | 11:10 am<br>*California* | **NCBI BLAST: Enhanced Web Usability through New Result Page and Effective Genomic Data Access**<br>*Digital Tools and Resources Session 3* |
| Wednesday | 11:50 am<br>*California* | **Federated Cloud Access to Datasets through Indexing and/or Graphs!**<br>*Digital Tools and Resources Session 4* |

# NCBI Genome Resources Workshop

Monday January 13, 2020, 12:50 – 3:00 pm, Pacific Salon 1

| Time | Topic |
|---|---|
| 12:55 – 1:15 | NCBI Wants Your Sequence Data! How Do I Get It There?<br>*Ilene Mizrachi* |
| 1:15 – 1:35 | Annotation of Eukaryote Genomes at NCBI<br>*Jinna Hoffman* |
| 1:35 – 1:55 | Accessing Homologous Gene Datasets at NCBI<br>*Nuala O'Leary* |
| 1:55 – 2:15 | The New PubMed Is Here!<br>*Kathi Canese* |
| 2:15 – 2:35 | Taxonomy Lookup; Data Retrieval: How to Find and Stream Genomic Data in the Cloud!<br>*Ben Busby* |

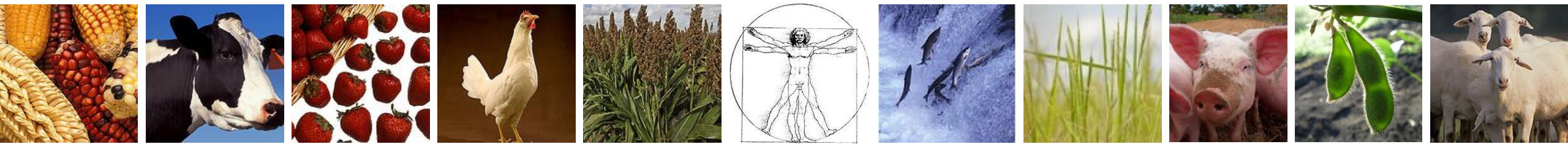Visit NCBI Booth **321**          Contact us **info@ncbi.nlm.nih.gov**

**NIH** U.S. National Library of Medicine
National Center for Biotechnology Information

Watch NCBI News for updates!
http://www.ncbi.nlm.nih.gov/news/
https://www.youtube.com/user/NCBINLM

# NCBI RefSeq Resources for Plant Genomics

Anjana R Vatsan(raina@nih.gov)

Functional Genomics Workshop

PAG XXVIII January 13, 2020

U.S. National Library of Medicine
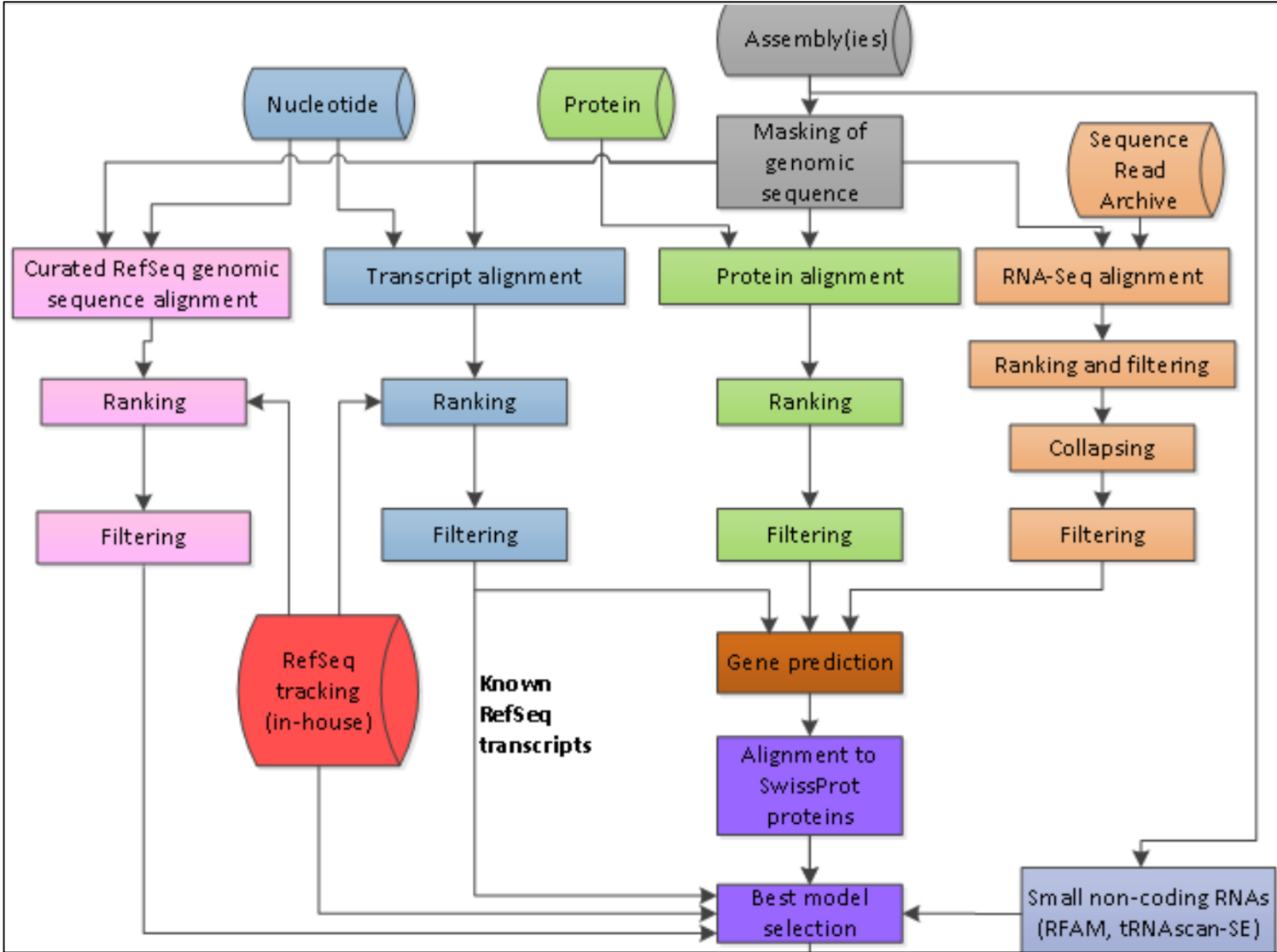*National Center for Biotechnology Information*

# RefSeq: NCBI Reference Sequence Database

RefSeq is a public database of nucleotide and protein sequences which are derived, in most part, from genome assemblies that are submitted to International Nucleotide Sequence Database Collaboration (INSDC), by one of the following methods.

- Computationally using the Genome Annotation Pipeline
- Manual curation
- Propagation from model organism databases
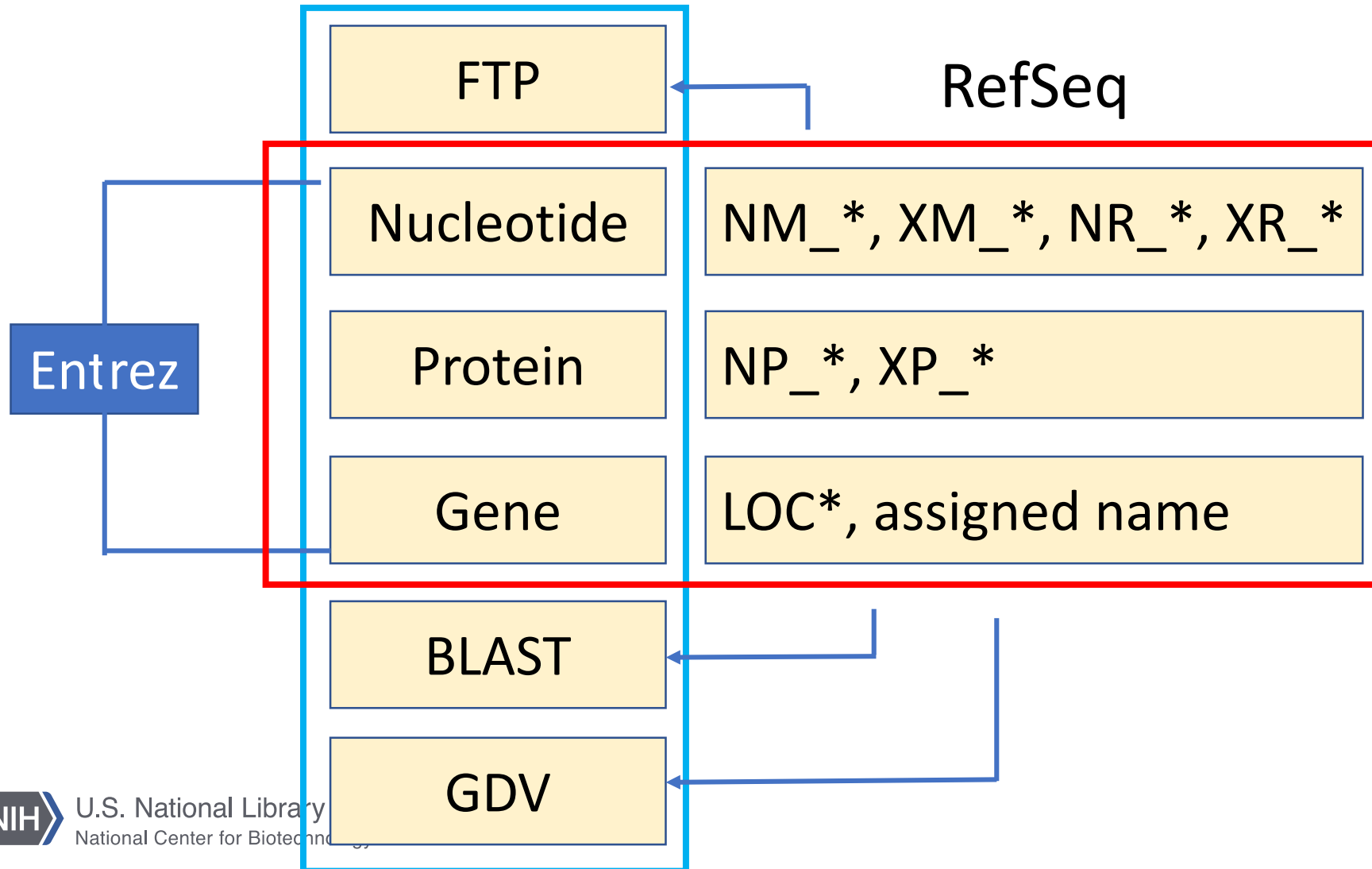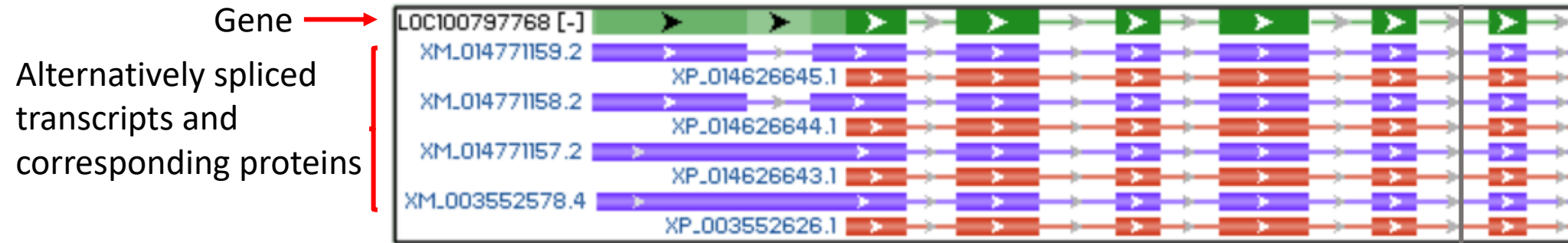    e.g. Arabidopsis thaliana

**https://www.ncbi.nlm.nih.gov/refseq/**

# Annotation Pipeline

# Annotation Output

**ftp://ftp.ncbi.nlm.nih.gov/genomes/refseq/plant/**



RefSeq

| | |
|---|---|
| FTP | |
| Nucleotide | NM_*, XM_*, NR_*, XR_* |
| Protein | NP_*, XP_* |
| Gene | LOC*, assigned name |
| BLAST | |
| GDV | |

Entrez

# RefSeq Data organization



Gene →

Alternatively spliced transcripts and corresponding proteins

LOC100797768 [-]
XM_014771159.2
XP_014626645.1
XM_014771158.2
XP_014626644.1
XM_014771157.2
XP_014626643.1
XM_003552578.4
XP_003552626.1

# 108 plant genomes have been annotated at NCBI



Monocots

Dicots

https://www.ncbi.nlm.nih.gov/taxonomy
https://www.ncbi.nlm.nih.gov/tools/treeviewer/

# Annotation of various organism groups

Assembly | soybean[orgn]

Best representative assembly chosen for annotation

Access assembly meta-data, statistical reports, and links to genomic sequence data

1. Glycine_max_v2.1
Organism: **Glycine max** (soybean)
Infraspecific name: Cultivar: Williams 82
Submitter: US DOE Joint Genome Institute (JGI-PGF)
Date: 2018/07/24
Assembly level: Chromosome
Genome representation: full
RefSeq category: representative genome
GenBank assembly accession: GCA_000004515.4 (**latest**)
RefSeq assembly accession: GCF_000004515.5 (**latest**)
IDs: 1832791 [UID] 6943708 [GenBank] 7001488 [RefSeq]

2. **Glycine max_Enrei_2.0**
Organism: **Glycine max** (soybean)
Infraspecific name: Cultivar: ENREI
Submitter: National Institute of Agrobiological Sciences
Date: 2015/08/06
Assembly level: Contig
Genome representation: full
GenBank assembly accession: GCA_001269945.2 (**latest**)
RefSeq assembly accession: n/a
IDs: 474001 [UID] 2242528 [GenBank]

3. glyma.Lee.gnm1
Organism: **Glycine max** (soybean)
Infraspecific name: Cultivar: Lee
Submitter: **Glycine max** cv Lee and Glycine soja PI 483463 sequencing consortium
Date: 2018/01/30
Assembly level: Chromosome
Genome representation: full
GenBank assembly accession: GCA_002905335.2 (**latest**)
RefSeq assembly accession: n/a
IDs: 2580391 [UID] 9338298 [GenBank]

4. Gmax_ZH13
Organism: **Glycine max** (soybean)
Infraspecific name: Cultivar: Zhonghuang 13
Submitter: Institute of Genetics and Developmental Biology , Chinese Academy of Science
Date: 2018/08/10

# Annotation Home Page

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Eutrema salsugineum (saltwater cress) | Eutsalg1_0 (GCF_000478725.1) | 100 | 2018-02-22 | 2018-02-26 | FTP | B | AR | GDV |
| Fragaria vesca (wild strawberry) | FraVesHawaii_1.0 (GCF_000184155.1) | 101 | 2015-03-02 | 2015-03-04 | FTP | B | AR | GDV |
| Glycine max (soybean) | Glycine_max_v2.1 (GCF_000004515.5) | 103 | 2018-08-07 | 2018-09-06 | FTP | B | AR | GDV |
| Glycine soja (wild soybean) | ASM419377v2 (GCF_004193775.1) | 100 | 2019-03-08 | 2019-03-12 | FTP | B | AR | GDV |
| Gossypium arboreum (tree cotton) | Gossypium_arboreum_v1.0 (GCF_000612285.1) | 100 | 2016-08-01 | 2016-08-11 | FTP | B | AR | GDV |
| Gossypium hirsutum (cotton) | ASM98774v1 (GCF_000987745.1) | 100 | 2016-05-09 | 2016-05-18 | FTP | B | AR | GDV |
| Gossypium raimondii (eudicots) | Graimondii2_0 (GCF_000327365.1) | 100 | 2015-04-29 | 2015-05-22 | FTP | B | AR | GDV |
| Helianthus annuus (common sunflower) | HanXRQr1.0 (GCF_002127325.1) | 100 | 2017-07-28 | 2017-08-07 | FTP | B | AR | GDV |
| Herrania umbratica (eudicots) | ASM216827v2 (GCF_002168275.1) | 100 | 2017-06-08 | 2017-06-09 | FTP | B | AR | GDV |
| Hevea brasiliensis (rubber tree) | ASM165405v1 (GCF_001654055.1) | 100 | 2017-07-16 | 2017-07-19 | FTP | B | AR | GDV |
| Ipomoea nil (Japanese morning glory) | Asagao_1.1 (GCF_001879475.1) | 100 | 2016-11-25 | 2016-11-29 | FTP | B | AR | GDV |
| Ipomoea triloba (trilobed morning glory) | ASM357664v1 (GCF_003576645.1) | 100 | 2019-10-11 | 2019-10-17 | FTP | B | AR | GDV |
| Jatropha curcas (eudicots) | JatCur_1.0 (GCF_000696525.1) | 101 | 2017-03-31 | 2017-04-06 | FTP | B | AR | GDV |

# Annotation Report

| Feature | Glycine_max_v2.1 |
|---|---|
| Genes and pseudogenes ⓘ | 59,906 |
| protein-coding | 46,993 |
| non-coding | 7,881 |
| transcribed pseudogenes | 376 |
| non-transcribed pseudogenes | 4,656 |
| genes with variants | 13,958 |
| immunoglobulin/T-cell receptor gene segments | 0 |
| other | 0 |
| mRNAs | 71,048 |
| fully-supported | 64,943 |
| with > 5% ab initio ⓘ | 5,225 |
| partial | 410 |
| with filled gap(s) ⓘ | 37 |
| known RefSeq (NM_) ⓘ | 7,593 |
| model RefSeq (XM_) | 63,455 |
| non-coding RNAs ⓘ | 14,358 |
| fully-supported | 11,288 |
| with > 5% ab initio ⓘ | 0 |
| partial | 0 |

### Detailed reports

| Feature | Count | Mean length (bp) | Median length (bp) | Min length (bp) | Max length (bp) |
|---|---|---|---|---|---|
| Genes | 54,874 | 4,017 | 2,958 | 54 | 287,040 |
| All transcripts | 85,406 | 1,872 | 1,655 | 18 | 16,833 |
| mRNA | 71,048 | 1,953 | 1,712 | 148 | 16,833 |
| misc_RNA | 4,407 | 2,389 | 2,080 | 197 | 8,834 |
| miRNA | 614 | 22 | 21 | 18 | 25 |
| tRNA | 752 | 74 | 73 | 71 | 93 |
| lncRNA | 6,360 | 1,521 | 1,146 | 54 | 11,067 |
| snoRNA | 1,687 | 105 | 107 | 64 | 228 |
| snRNA | 136 | 151 | 158 | 100 | 198 |
| rRNA | 402 | 1,615 | 1,807 | 104 | 3,470 |
| Single-exon transcripts ⓘ | 7,813 | 1,274 | 1,044 | 148 | 8,719 |
| coding transcripts (NM_/XM_) | 7,780 | 1,275 | 1,046 | 148 | 8,719 |

https://www.ncbi.nlm.nih.gov/genome/annotation_euk/all

# Annotation Output

**ftp://ftp.ncbi.nlm.nih.gov/genomes/refseq/plant/**

RefSeq

| | |
|---|---|
| FTP | |
| Nucleotide | NM_*, XM_*, NR_*, XR_* |
| Protein | NP_*, XP_* |
| Gene | LOC*, assigned name |
| BLAST | |
| GDV | |

Entrez

# Data Access-Entrez search

https://www.ncbi.nlm.nih.gov/search/?term=   150%   Search

# NIH  U.S. National Library of Medicine
National Center for Biotechnology Information

Log in

## Search NCBI

Soybean F ✕   **Search**

Glycine max FAD2-1B

**Glycine max FT2A**

Glycine max FAD2-2

Glycine max FG3

Glycine max FLS1

Glycine max FAD2-1A

Glycine max FAD3B

Glycine max FNSII-1

# News

## Recent blog posts

**DECEMBER 23, 2019**
Mitochondrial COX1 submission improv
now live in submission portal!

**DECEMBER 20, 2019**
ClinVar Celebrates 1 Million Submissions

**DECEMBER 18, 2019**
BLAST+ 2.10.0 now available with improved
composition-based statistics

**DECEMBER 17, 2019**
Genome Workbench is now in the cloud!

the Past Faster

he what ancient
people ate and monitor historical sites from the sky

**NPR News**   DECEMBER 25, 2019

## A Young Mississippi Woman's Journey Through A Pioneering Gene-Editing Experiment

NPR tells the exclusive, behind-the-scenes story of the first person with a genetic disorder to be treated in the United ...

# NCBI Databases

## Literature
The World's largest repository of medical and

## Genes
Gene sequences and annotations used as references

## Proteins
Protein sequences, 3-D structures, and tools for the

GENE

Was this helpful? 👍 👎

# Results by database

Results found in 10 databases

| Literature | | Genes | | Proteins | |
|---|---|---|---|---|---|
| Bookshelf | 0 | Gene | 1 | Conserved Domains | 0 |
| MeSH | 0 | GEO DataSets | 0 | Identical Protein Grou | 14 |
| NLM Catalog | 0 | GEO Profiles | 0 | Protein | 166 |
| PubMed | 10 | HomoloGene | 0 | Protein Clusters | 0 |
| PubMed Central | 34 | PopSet | 1 | Sparcle | 0 |
| | | | | Structure | 0 |

NIH

Was this helpful? 👍 👎

# FT2A – protein FLOWERING LOCUS T

Glycine max (soybean)

Also known as: GLYMA_16G150700, E9, FT, FT3, FTL3, GmFT2a

GeneID: 100814951

RefSeq transcripts (2)    RefSeq proteins (2)    PubMed (10)

| Genome Browser | BLAST | Download |

---

## RefSeq Sequences    —

| NM_001253256.2 | 899 | NP_001240185.1 | 176 | |
| XM_006598696.3 | 920 | XP_006598759.1 | 158 | X1 |

---

# Results by database

# Data access – Gene and RefSeq

## Bibliography

### Related articles in PubMed

1. CRISPR/Cas9-mediated targeted muta...
   Cai Y, *et al.* Plant Biotechnol J, 2018 Jan.
2. GmFT2a polymorphism and maturity ...
   Jiang B, *et al.* PLoS One, 2013. PMID 241...
3. GmFT2a, a soybean homolog of FLOW...
   Sun H, *et al.* PLoS One, 2011. PMID 2219...
4. Mutagenesis of GmFT2a and GmFT5a...
   Cai Y, *et al.* Plant Biotechnol J, 2020 Jan.
5. Functional divergence between soybea...
   Takeshima R, *et al.* J Exp Bot, 2019 Aug 7.

See all (10) citations in PubMed

### GeneRIFs: Gene References Into Function

What's a GeneRIF?

1. GmFT1a expression was induced by l...
   expression of flowering promoters Gm
2. Expression of the FLOWERING LOCU
3. Although GmFT2a is a key flowering g
4. GmFT2a expression is associated with

**Submit:** New GeneRIF   Correction

---

GeneRIF (Gene Reference Into Function) enables interested scientists to enrich the functional
annota...
and th...
annota...

Please...
will no...

For ad...

To sug...
please...

Each G...
GeneR

Submi...

Gene I
10081
PubMe

GeneR

---

## NCBI — Feedback for Gene and Reference Sequences (RefSeq)

Make suggestions, submit additions and corrections, or ask for help concerning Gene or Reference Sequence (RefSeq) records. See additional information: Gene Home Page, RefSeq Home Page.

Do not use this form to report a problem in PubMed or GenBank. Do not use this form to submit sequence data to NCBI.

Additional contacts:

- PubMed - report typographical or other errors in citations
- GenBank submission documentation - submit sequence data to GenBank
- NCBI help - contact us about other NCBI resources

### What would you like to do?

Return to previous page

Add a GeneRIF - Add a publication with a functional comment to a Gene record.

**Additions**

○ Report a new gene that is not yet available in Gene
○ Request addition of a RefSeq transcript, protein, or pseudogene record
○ Contribute a summary describing the function of the gene

**Corrections**

○ Correct or update a Gene record (please provide the GeneID)
○ Correct or update a RefSeq record (please provide the accession.version)
○ Report a publication that is incorrectly associated with a Gene or RefSeq (please provide the PubMed ID

# Data access – Gene and RefSeq



**FT2A**   protein FLOWERING LOCUS T [ *Glycine max* (soybean) ]

Gene ID: 100814951, updated on 14-Dec-2019

### Summary

| | |
|---|---|
| Gene symbol | FT2A |
| Gene description | protein FLOWERING LOCUS T |
| Locus tag | GLYMA_16G150700 |
| Gene type | protein coding |
| RefSeq status | VALIDATED |
| Organism | Glycine max |
| Lineage | Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta; Spermatophyta; Magnoliopsida; eudicotyledons; Gunneridae; Pentapetalae; rosids; fabids; Fabales; Fabaceae; Papilionoideae; 50 kb inversion clade; NPAAA clade; indigoferoid/millettioid clade; Phaseoleae; Glycine; Soja |
| Also known as | E9; FT; FT3; FTL3; GmFT2a |

### Genomic context

Location:   chromosome: 16                                             See FT2A in Genome Data Viewer

Exon count:   4

| Annotation release | Status | Assembly | Chr | Location |
|---|---|---|---|---|
| 103 | current | Glycine_max_v2.1 (GCF_000004515.5) | 16 | NC_038252.1 (31109897..31114974) |
| 102 | previous assembly | Glycine_max_v2.0 (GCF_000004515.4) | 16 | NC_016103.2 (31109907..31114981) |
| 101 | previous assembly | V1.1 (GCF_000004515.3) | 16 | NC_016103.1 (30741587..30746627) |

**Chromosome 16 - NC_038252.1**

[ 31082789 ▶          [ 31151913 ▶

LOC100814418          LOC102664067                    FT2B
     LOC100305576        FT2A
              LOC100787142

### Genomic regions, transcripts, and products

### Bibliography

### Variation

### Pathways from PubChem

### General gene information

### General protein information

### NCBI Reference Sequences (RefSeq)

### Related sequences

**Table of contents**
Summary
...oducts
General gene information
    Homology
General protein information
NCBI Reference Sequences (RefSeq)
Related sequences
Additional links

**Genome Browsers**
Genome Data Viewer

**Related information**
BioProjects
BioSystems
Conserved Domains
Full text in PMC
Full text in PMC_nucleotide
Gene neighbors
Genome
Nucleotide
Protein
PubMed
PubMed (GeneRIF)
PubMed(nucleotide/PMC)
RefSeq Proteins
RefSeq RNAs
SNP: GeneView
Taxonomy
UniGene

**GLYMA_16G150700**

**Gene display includes reference sequences and various links and tools for the study of gene expression and function.**

**Location of gene on the chromosome**
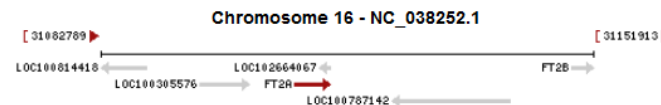
FT2A has two transcript variants

NM_001253256.2
XM_006598696.3

# Difference between N* and X* accessions

- N* accessions are used for **known RefSeqs**. This category is supported by manual curation. Records are primarily derived from INSDC cDNAs, EST, and Transcript Shotgun Assembly (TSA) records.

- X* accessions are **predicted models**. This category is computationally predicted based on aligned evidence. Records are primarily derived from genomic sequence. The vast majority of XMs are fully supported by experimental evidence, and for most species **they are on par, quality-wise, with the NMs.**

## NM_001253256.2

```
COMMENT        VALIDATED REFSEQ: This record has undergone validation or
               preliminary review. The reference sequence was derived from
               ACUP03010190.1 and KJ607992.1.
               On Jul 19, 2017 this sequence version replaced NM_001253256.1.

               ##Evidence-Data-START##
               Transcript exon combination :: EU287455.1, AB550122.1 [ECO:000033
               RNAseq introns              :: single sample supports all introns
                                              SAMN02009287, SAMN02215336
                                              [ECO:0000348]
               ##Evidence-Data-END##
PRIMARY        REFSEQ_SPAN         PRIMARY_IDENTIFIER PRIMARY_SPAN        COMP
               1-77                ACUP03010190.1     61121-61197
               78-608              KJ607992.1         1-531
               609-899             ACUP03010190.1     65802-66092
FEATURES             Location/Qualifiers
     source          1..899
                     /organism="Glycine max"
                     /mol_type="mRNA"
                     /cultivar="Williams 82"
                     /db_xref="taxon:3847"
                     /chromosome="16"
```

## XM_006598696.3

```
COMMENT        MODEL REFSEQ: This record is predicted by automated computational
               analysis. This record is derived from a genomic sequence
               (NC_038252.1) annotated using gene prediction method: Gnomon.
               Also see:
                   Documentation of NCBI's Annotation Process

               On Aug 15, 2018 this sequence version replaced XM_006598696.2.

               ##Genome-Annotation-Data-START##
               Annotation Provider         :: NCBI
               Annotation Status           :: Full annotation
               Annotation Name             :: Glycine max Annotation Release 103
               Annotation Version          :: 103
               Annotation Pipeline         :: NCBI eukaryotic genome annotation
                                              pipeline
               Annotation Software Version :: 8.1
               Annotation Method           :: Best-placed RefSeq; Gnomon
               Features Annotated          :: Gene; mRNA; CDS; ncRNA
               ##Genome-Annotation-Data-END##
FEATURES             Location/Qualifiers
     source          1..920
                     /organism="Glycine max"
                     /mol_type="mRNA"
                     /cultivar="Williams 82"
                     /db_xref="taxon:3847"
                     /chromosome="16"
                     /tissue_type="callus"
```

# Data curation—how do we maintain the quality of our data

PRR3A gene (GeneID: **100785796)**

Low Quality protein, corrected it based on PMID: 30418611*



NM_001377264.1
XM_014773179.1

*Li et. al. Plant Cell Physiol. 2019 Feb 1;60(2):407-420. Characterization of Two Growth Period QTLs Reveals Modification of PRR3 Genes During Soybean Domestication.

```
LOCUS       XM_014773179            3829 bp    mRNA    linear   PLN 31-AUG-2018
DEFINITION  PREDICTED: Glycine max two-component response regulator-like PRR37
            (LOC100785796), mRNA.
ACCESSION   XM_014773179
VERSION     XM_014773179.1
DBLINK      BioProject: PRJNA48389
KEYWORDS    RefSeq; corrected model.
SOURCE      Glycine max (soybean)
```

```
        2526-3131          ACUP03013092.1      75272-75877        c
        3132-3132          "N"                 1-1
        3133-3284          ACUP03013092.1      75120-75271        c
        3285-3358          ACUP03013092.1      74576-74649        c
        3359-3829          ACUP03013092.1      74017-74487        c
FEATURES             Location/Qualifiers
     source          1..3829
                     /organism="Glycine max"
                     /mol_type="mRNA"
                     /cultivar="Williams 82"
                     /db_xref="taxon:3847"
                     /chromosome="Unknown"
                     /tissue_type="callus"
     gene            1..3829
                     /gene="LOC100785796"
                     /note="The sequence of the
                     modified relative to its s
                     represent the inferred CDS
                     Derived by automated compu
                     prediction method: Gnomon.
                     similarity to: 31 ESTs, 2
                     the annotated genomic feat
                     including 147 samples with
                     introns"
                     /db_xref="GeneID:10078579
     CDS             1086..3437
                     /gene="LOC100785796"
                     /note="The sequence of the model RefSeq protein was
                     modified relative to its source ge
                     represent the inferred CDS: insert
                     /codon_start=1
                     /product="LOW QUALITY PROTEIN: two-component response
```
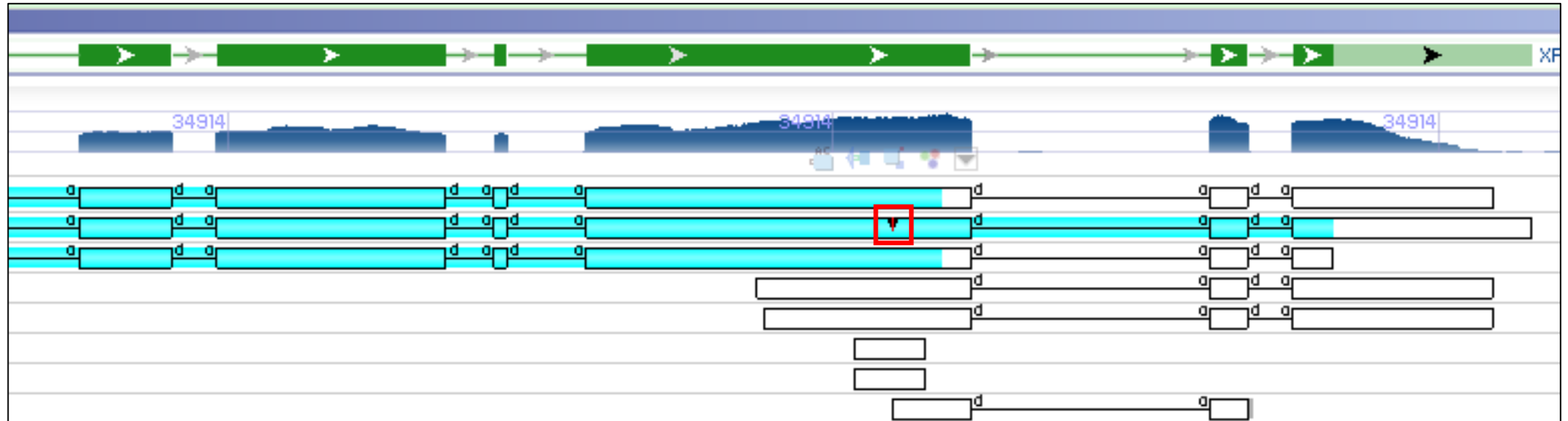
3132-3132          "N"          1-1

The sequence of the model RefSeq **transcript was modified** relative to its source genomic sequence to represent the inferred CDS: **inserted 2 bases in 2 codon**; Derived by automated computational analysis using gene prediction method: Gnomon. Supporting evidence includes similarity to: 31 ESTs, 2 Proteins, and 100% coverage of the annotated genomic feature by RNASeq alignments.

/product="LOW QUALITY PROTEIN: two-component response

# Data curation: Add value based on publications
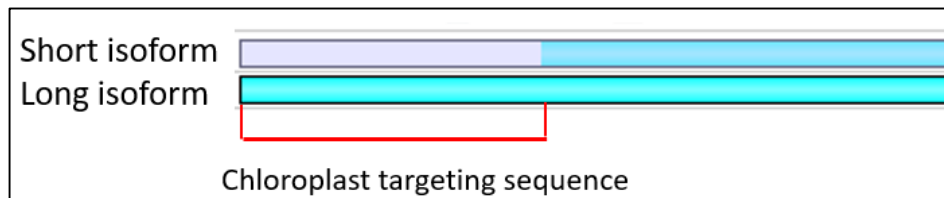
Adding value to data: Updated HPPD gene (GeneID:100101901) to create two isoforms, short and long, based on PMID: 25192697* .

Updating locus_type: Updated LUXa gene from coding to non-coding based on PMID:28878247**

*Seihl et. al. Plant Physiol. 2014 Nov;166(3):1162-1176. Broad 4-hydroxyphenylpyruvate dioxygenase inhibitor herbicide tolerance in soybean with an optimized enzyme and expression cassette
**Liew et. Al. Sci Rep. 2017 Sep 6;7(1):10605. A novel role of the soybean clock gene LUX ARRHYTHMO in male reproductive development.

U.S. National Library of Medicine
National Center for Biotechnology Information

# Data Analysis using BLAST

ICA2 gene (AT2G32320) involved in growth and flowering time plasticity in relation to temperature in Arabidopsis; PMID: 30992321*

Gene | Gene | At2g32

Create

**AT2G32320  tRNAHis guanylyltransfer**

Gene ID: 817793, updated on 25-Oct-2019

▼ **Summary**

▼ **Genomic context**

▼ **Genomic regions, transcripts, and products**

▼ **Bibliography**

▼ **Variation**

▼ **Pathways from PubChem**

▼ **General gene information**

▼ **General protein information**

▲ **NCBI Reference Sequences (RefSeq)**   ⊗ ?

1. NM_128791.4 → NP_180791.3  **tRNAHis guanylyltransferase [Arabidopsis thaliana]**

   See identical proteins and their annotated locations for NP_180791.3

   **Status: REVIEWED**

   UniProtKB/Swiss-Prot   F4ISV6
   Conserved Domains (2) summary

   | pfam04446 | Thg1; tRNAHis guanylyltransferase |
   Location:275 → 401 |

   | pfam14413 | Thg1C; Thg1 C terminal domain |
   Location:405 → 510 |

2. NM_001161072.1 → NP_001154544.1  **tRNAHis guanylyltransferase [Arabidopsis thaliana]**

   See identical proteins and their annotated locations for NP_001154544.1

   **Status: REVIEWED**

   UniProtKB/Swiss-Prot   F4ISV6
   Conserved Domains (3) summary

   | COG4021 | Thg1; tRNA(His) 5'-end guanylyltransferase [Translation, ribosomal structure and biogenesis] |
   Location:286 → 531 |

   | pfam04446 | Thg1; tRNAHis guanylyltransferase |
   Location:287 → 413 |

   | pfam14413 | Thg1C; Thg1 C terminal domain |
   Location:417 → 522 |

3. NM_001161073.1 → NP_001154545.1  **tRNAHis guanylyltransferase [Arabidopsis thaliana]**

*Mendez-Vigo et. al. Plant Cell. 2019 Jun;31(6):1222-1237. Genetic Interactions and Molecular Evolution of the Duplicated Genes *ICARUS2* and *ICARUS1* Help Arabidopsis Plants Adapt to Different Ambient Temperatures.

## Enter Query Sequence

**Enter accession number(s), gi(s), or FASTA sequence(s)** ⑦     Clear          **Query subrange** ⑦

```
NP_180791.3
```

From [          ]

To [          ]

**Or, upload file**     [ Browse... ]  No file selected.     ⑦

**Job Title**     [ NP_180791:tRNAHis guanylyltransferase [Arabidopsis... ]

Enter a descriptive title for your BLAST search ⑦

☐ **Align two or more sequences** ⑦

## BLAST results will be displayed in a new format by default

**New**

You can always switch back to the Traditional Results page.

## Choose Search Set

**Database**     [ Non-redundant protein sequences (nr)            ∨ ] ⑦

**Organism**
*Optional*     [ soybean (taxid:3847)                    ]  ☐ exclude  [ + ]

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown. ⑦

# Accessing RefSeq from BLAST output

Select RefSeqs

| | | | | | | |
|---|---|---|---|---|---|---|
| ☑ | tRNA(His) guanylyltransferase 2 [Glycine max] | 643 | 643 | 99% | 0.0 | 57.63% | XP_003536324.1 |
| ☑ | tRNA(His) guanylyltransferase 2 [Glycine max] | 637 | 637 | 99% | 0.0 | 57.06% | XP_003547793.1 |

| Expect | Method | Identities | Positives | Gaps |
|---|---|---|---|---|
| 0.0 | Compositional matrix adjust. | 302/524(58%) | 395/524(75%) | 17/524(3%) |

**Information**

Gene - associated gene details

Genome Data Viewer aligned genomic context

Identical Proteins - Identical proteins to XP_003536324.1

```
MANSKYEYVKSFEVEDEVMFPNLIIIRIDGRDFSRFSQVHKFEKPNDETSLNL
MANSKYEYVK FEVEDE MFPN+I++ I        +   K  KP+D  +L L
MANSKYEYVKCFEVEDEAMFPNIILVWI---------KASKLHKPHDSNTLKL

VLVEYPDIVFAYGYSDEYSFVFKKASRFYQRRASKILSLVASFFAAVYVTKWK
VL EY D+VFAYG+SDEY+FVFKK S+F++RRASK+LS++ SFF++V+V KW
VLEEYADVFAYGFSDEYTFVFKKTSKFHERRASKVLSIITSFFSSVFVRKWD

LEYAPSFASKVVSCASVEVLQAYLAWRQHDCHISNQYDTCLWMLVKSGKTLSE
L+  PS  +V++CAS++ LQAYL WRQ  CH+SNQ++ CLW LV+ G   E
LQCHPSLHGRVIACASIKALQAYLLWRQTICHLSNQHEQCLWRLVERGMNEKE

TQKQQRNELLFQQFGINYKMLPVLFRQGSCLFKTKLEETVKHDENGKPVKRLR
 +K   N LLF +F +NY  L  + RQGSC+ KT  E+TVK+ +NG P+KR R
FEKSDLNNLLFDEFNVNYNTLEPILRQGSCVLKTTGEDTVKYTDNGAPIKRHR
```

| | | |
|---|---|---|
| ☑ | hypothetical pr | RH07905.1 |
| ☑ | hypothetical pr | RH08092.1 |
| ☑ | tRNA(His) gua | P_025981896.1 |
| ☑ | hypothetical pr | RH64848.1 |
| ☑ | hypothetical pr | RH06105.1 |
| ☑ | hypothetical pr | RH06104.1 |
| ☑ | midasin [Glycin | P_014624284.1 |
| ☑ | midasin [Glycin | P_006583141.1 |

NC_038246.1

Select an assembly to change view

▼ Ideogram View

Unplaced/unlocalized scaffolds: 1,170

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

16 17 18 19 20 Pltd MT

▼ Search

🔍 Location, gene or phenotype

Enter a location, gene name or phenotype

▶ Search examples:

▶ User Data and Track Hubs

▶ BLAST

▶ Add Tracks

▶ History

🔧 Tools ▾   ⚙ Tracks ▾   ⬇ Download ▾   ↻   ❓ ▾

150 K  151 K  152 K  153 K  154 K  155 K  156 K  157 K  158 K  159 K  160 K  161 K  162 K  163 K

Genes, NCBI Glycine max Annotation Release 103, 2018-09-10

LOC100804661                                LOC100802539

XM_014762856.2                     XM_003536276.4                    XM_003536276.4
XP_014618342.1                                                        XP_003536324.1

PPR_2   PPR   CLH                                                     Thg1
PPR_2        PPR_2                                                    Thg1      Thg1
PPR_2   PPR_2   PPR                                                   Thg1
PPR   PPR_2   PPR   PPR                                               Thg1C
PPR   PPR                                                             Thg1C
PPR        PPR                                                        XM_006588467.3
PPR repeat   PPR repeat                                               XP_006588530.1
PPR repeat   PPR repeat
PPR repeat   PPR repeat                                               Thg1
PPR repeat   PPR repeat                                               Thg1   Thg1
PPR repeat   PPR repeat
PPR repeat   PPR repeat                          Thg1
PPR repeat                                region: Thg1
PPR repeat                                Comment: tRNA(His) 5'-end guanylyltransferase [Translation, ribosomal structure and biogenesis]
PPR repeat                                Location: 2..227
PPR repeat              XM_006588468.3     Length: 226 aa
PPR repeat
PPR repeat                                Links & Tools
                                          View CDD: 226508

Reformat    Format: Hypertext    Row Display: up to 10    Color Bits: 2.0 bit    Type Selection: the most diverse members

3OTB_A     138 QTLKDYLSWRQADCHINNLYNTVFWALIQQSGLTPVQAQGRLQGTLAADKNEILFSEFNINYNNELPMYRKGTVLIWqk- 216 human
Q9Y7T3     137 KVLRDYLHWRQVDCHINNLYNTTFWMLILKGGFTNTQAEEYLKGTVSAEKHEILFSKFGINYNFEPEIYKKGSIWIRep- 215 Schizosaccharom...
XP_002175352 137 ssLRDYLSWRQADCHINNLYNTTFWALRLQGKMSNREAEERLKGTVSADKHEILFSQFGINYNNEPEMYKKGTIFTRkpa 216 Schizosaccharom...
XP_004366437 138 QNMRDYLSWRQADTHINNMYNTCYWALVLQGGCTPKEAEQTLCGTLSDAKNEILFTRFNINYNNLPQMYRKGSVIYRkm- 216 Dictyostelium f...
XP_007402901 138 KEVRDYFAWRQADTHINNLYNTTFWALVQQGGQTTTEAHSTLRGTVSKQHEVLFSRFGINYNDIAERYRKGSVLVRek- 216 Phanerochaete c...
EPY49808   137 sVLRDYLNWRQVDCHINNLYNTTFWALIQKGGLTNTKAEEYLKGTISSQKHEILFSQFHINYNNEREIYKKGSIWVRep- 215 Schizosaccharom...
XP_002468831 138 KEIRDYFSWRQADTHINNLYNTIFWALVQQGGETTTQAHATLRGTVSGTKNEMLHSRFGINYNTIPARYRKGSVLVQer- 216 Postia placenta...
CBK22716   137 QNIRDYISWRQADTHINNLYNTCFWALVQRGNETTTSAEKILNGTLSSEKNEILFSRFGINYNNEPEVFKKGSIVIRet- 215 Blastocystis ho...
EMS20245   143 aEVRDYLRWRQVDTHINNMYNTVFWALVLQGGRTPVEAEQELSGTISSQKQEILFSQFGINYNNLEPMYRKGSLVIWee- 221 Rhodosporidium ...
EJT49864   167 KEIRDYFAWRQADTHINNLYNTCFWALV-KAGRTPREANKELQGTNSKDKNEMLFSEFGINYNDIDPFYRKGSVLVRidp 245 Trichosporon as...

3OTB_A     217 -----------------------------------vD-----------------Evmtk------------ 223 human
Q9Y7T3     216 -----------------------------------idQ-----------------E----------- 219 Schizosaccharom...
XP_002175352 217 -----------------------------------dgD-----------------D----------- 220 Schizosaccharom...
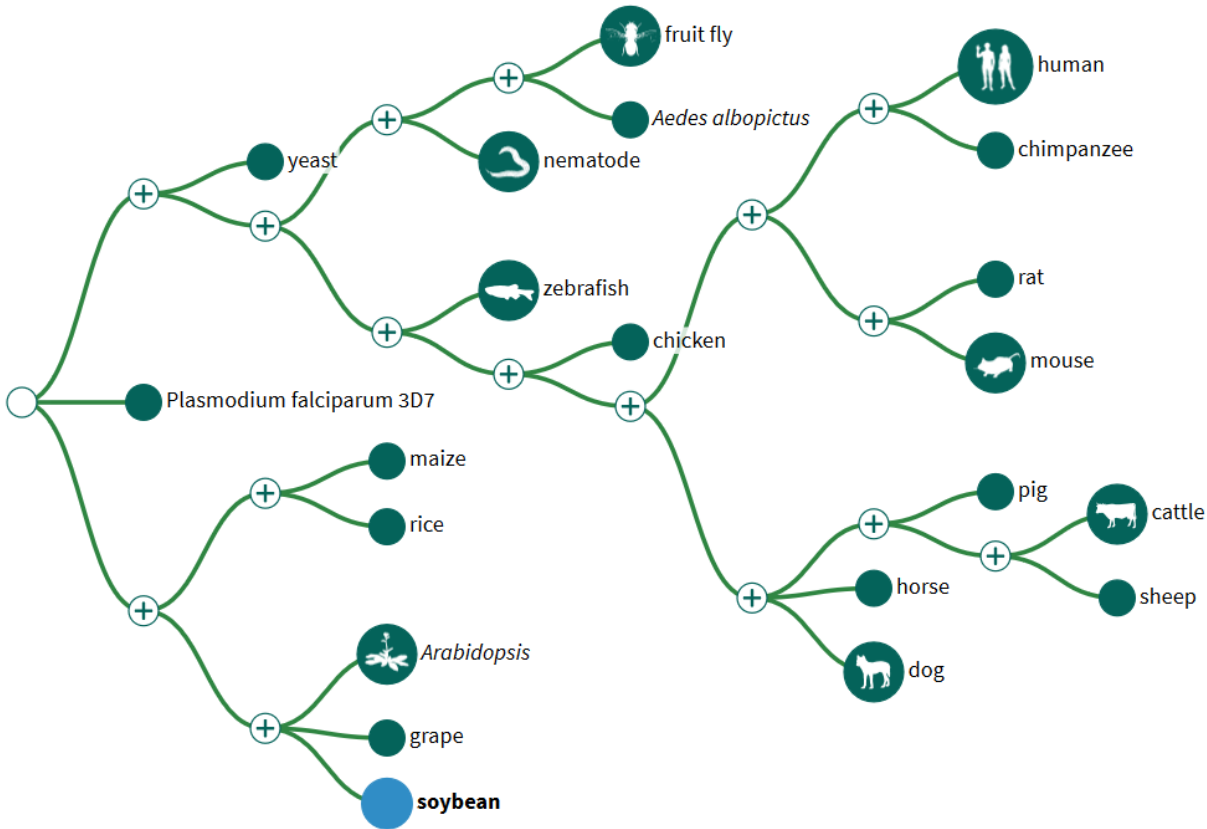XP_004366437 217 -----------------------------------vT-----------------E----------- 219 Dictyostelium f...

💬 Feedback

# Genome Data Viewer

GDV is a genome browser supporting the exploration and analysis of more than 740 eukaryotic RefSeq genome assemblies. ℹ

**Select organism**

Glycine max (soybean)



## *Glycine max* (soybean) genome

**Search in genome**

Flowering Locus 🔍

| Genes | Other |

| Name | Location |
|---|---|
| FT2A | Chr16: 31.11M - 31.11M |
| FT5A | Chr16: 4.136M - 4.138M |
| FT2C | Chr2: 6.099M - 6.111M |
| LOC100804540 | Chr5: 34.27M - 34.29M |
| FTL4 | Chr8: 47.46M - 47.46M |
| FT3A | Chr16: 4.162M - 4.165M |
| FT6 | Chr2: 47.47M - 47.47M |

Examples: KTI3, chr8:45734000-45738000, DNA repair

**Assembly**

Glycine_max_v2.1 ⌄

[ Browse genome ]  [ BLAST genome ]

### Assembly details

| | |
|---|---|
| **Name** | Glycine_max_v2.1 |
| **RefSeq accession** | GCF_000004515.5 |
| **GenBank accession** | GCA_000004515.4 |

💬 Feedback

# Genome Data Viewer

Select organism

Glycine max (soybean)



*Glycine max* (soybean) genome

Search in genome

Location, gene or phenotype

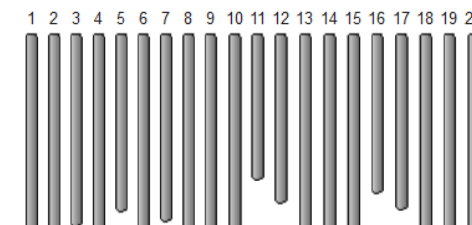Examples: KTI3, chr8:45734000-45738000, DNA repair

Assembly

Glycine_max_v2.1

Ca1 Ca2 Ca3 Ca4 Ca5 Ca6 Ca7 Ca8 Pltd

**Exon Navigator**

There are too many (3361) genes in the region. Please narrow the region to enable exon navigation.

NC_021160.1

|2 M |4 M |6 M |8 M |10 M |12 M |14 M |16 M |18 M |20 M |22 M |24 M |26 M |28 M |30 M |32 M |34 M

Genes, NCBI Cicer arietinum Annotation Release 102, 2018-12-14
ARF2        ELR19    ARF4      NAC4   CYP81E5      LOC101500381
HSFA4C              ARF5            CYP81E4

▼ Search

Location, Gene, Phenotype

▸ Search examples:

▸ User Data and Track Hubs

▼ BLAST                    ?

OCD73ME0014

▸ Add Tracks

▸ History

coverage, aggregate (filtered), NCBI Cicer arietinum Annotation Release 102 - log base 2 scaled

RNA-seq intron-spanning reads, aggregate (filtered), NCBI Cicer arietinum Annotation Release 102 - log base 2 scaled

RNA-seq intron features, aggregate (filtered), NCBI Cicer arietinum Annotation Release 102

|2 M |4 M |6 M |8 M |10 M |12 M |14 M |16 M |18 M |20 M |22 M |24 M |26 M |28 M |30 M |32 M |34 M

NC_021160.1: 1..48M (48,359,943 nt)

fruit fly

*Aedes albopictus*

human

Institute (JGI-

Release date          2018-09-06

1  2  3  4  5  6  7  8  9  10 11 12 13 14 15 16 17 18 19 20

grape

**soybean**

**NIH** U.S. National Library of Medicine
National Center for Biotechnology Information

📋 **Feedback**

# Genome Data Viewer

Select organism

Glycine max (soybean)





**NIH** U.S. National Library of Medicine
National Center for Biotechnology Information

---

*Glycine max* (soybean) genome

Search in genome

Location, gene or phenotype 🔍

Examples: KTI3, chr8:45734000-45738000, DNA repair

Assembly

Glycine_max_v2.1 ▼

**Browse genome**   **BLAST genome**

**Assembly details**

| | |
|---|---|
| **Name** | Glycine_max_v2.1 |
| **RefSeq accession** | GCF_000004515.5 |
| **GenBank accession** | GCA_000004515.4 |
| **Download via FTP** | RefSeq, GenBank |
| **Submitter** | US DOE Joint Genome Institute (JGI-PGF) |
| **Level** | Chromosome |
| **Category** | Representative genome |

**Annotation details**

| | |
|---|---|
| **Annotation Release** | 103 |
| **Release date** | 2018-09-06 |

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20

💬 Feedback

| Job Title | **ref\|NP_180791\|** |
|---|---|
| RID | 0CG217GJ014 *Search expires on 12-29 00:32 am* Download All ✔ |
| Program | TBLASTN ❓ Citation ✔ |
| Database | genomic/3847/GCF_000004515.5 See details ✔ |
| Query ID | NP_180791.3 |
| Description | tRNAHis guanylyltransferase [Arabidopsis tha ... |
| Molecule type | amino acid |
| Query Length | 525 |
| Other reports | ❓ |

## Filter Results

**Organism** *only top 20 will appear* ☐ exclude

Type common name, binomial, taxid or group name

✚ Add organism

| Percent Identity | E value | Query Coverage |
|---|---|---|
| [ ] to [ ] | [ ] to [ ] | [ ] to [ ] |

**Filter** **Reset**

---

**Descriptions** | Graphic Summary | Alignments | Taxonomy

### Sequences producing significant alignments

Download ✔   Manage Columns ✔   Show 100 ✔ ❓

☑ select all   *7 sequences selected*                GenBank   Graphics

| | Description | Max Score | Total Score | Query Cover | E value | Per. Ident | Accession |
|---|---|---|---|---|---|---|---|
| ☑ | Glycine max cultivar Williams 82 chromosome 10, Glycine_max_v2.1, whole genome shotgun sequence | 643 | 643 | 99% | 0.0 | 57.63% | NC_038246.1 |
| ☑ | Glycine max cultivar Williams 82 chromosome 16, Glycine_max_v2.1, whole genome shotgun sequence | 637 | 1844 | 99% | 0.0 | 57.06% | NC_038252.1 |
| ☑ | Glycine max cultivar Williams 82 chromosome 3, Glycine_max_v2.1, whole genome shotgun sequence | 207 | 366 | 72% | 4e-56 | 52.94% | NC_016090.3 |
| ☑ | Glycine max cultivar Williams 82 chromosome 13, Glycine_max_v2.1, whole genome shotgun sequence | 72.0 | 112 | 29% | 2e-11 | 40.00% | NC_038249.1 |
| ☑ | Glycine max cultivar Williams 82 chromosome 9, Glycine_max_v2.1, whole genome shotgun sequence | 45.8 | 81.2 | 9% | 0.002 | 62.96% | NC_038245.1 |

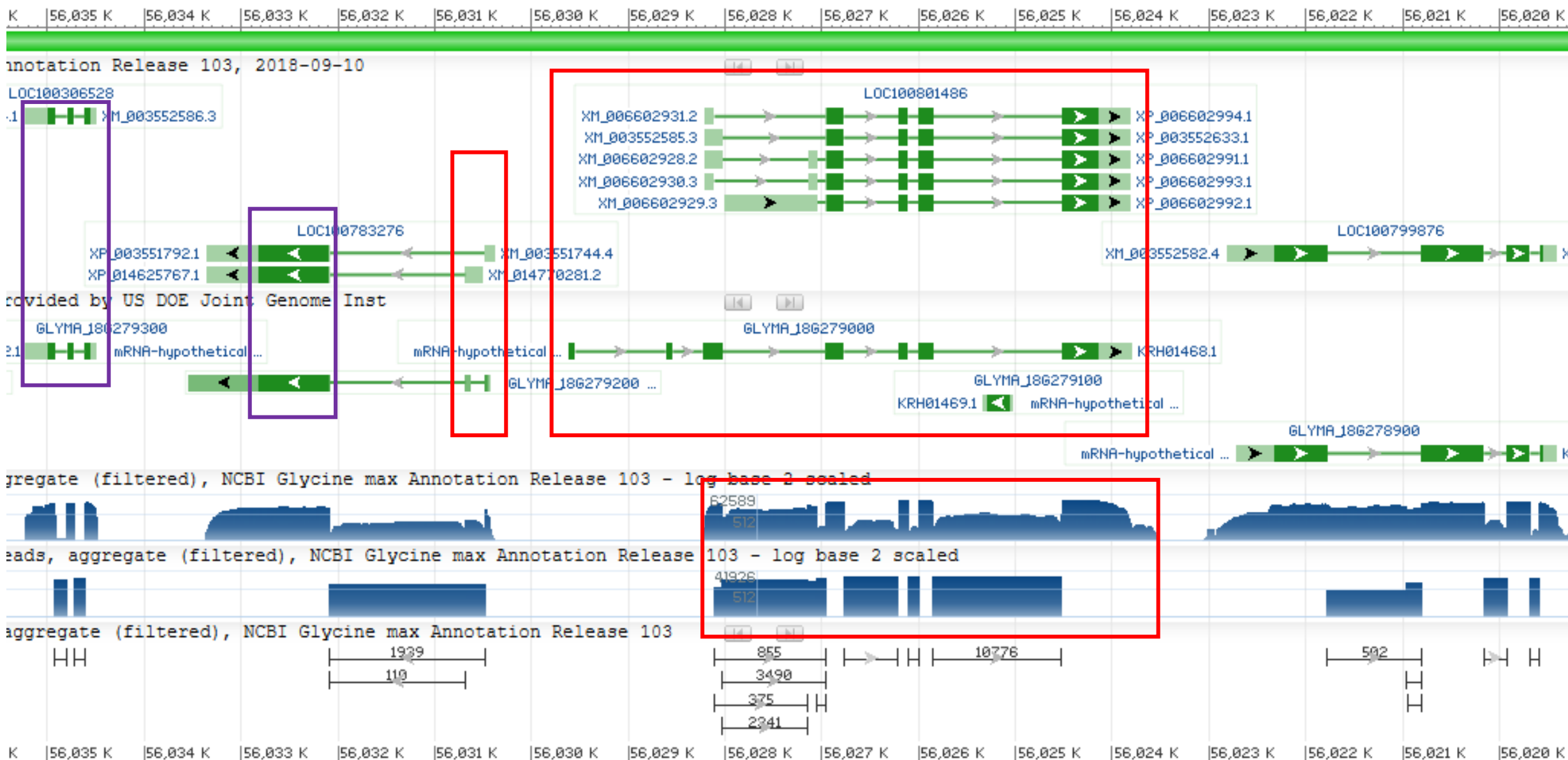Pick the exon

Pick the exon

Multiple hits

# Data Quality

GeneID: 100801486 YIF1B-like protein

*Li et. al., while assessing flavonoid biosynthesis pathway (FBP) genes in various Solanaceae species found that NCBI annotations were generally longer and concluded that

"Overall, the large majority of FBP homologs detected from these various annotation sources were in close agreement, but when they differed, **homologs from the NCBI annotations were generally longer and more abundant than from the genome-specific databases. These results suggest that reannotation of genome sequences using a unified annotation strategy, such as employed by the NCBI pipeline, may be preferable for improved consistency during comparative genomics research."**

*Li et. al. Genes (Basel). 2019 Jul 25;10(8). Assessing Anthocyanin Biosynthesis in Solanaceae as a Model Pathway for Secondary Metabolism.

# Thank you.

### RefSeq/Gene

**Terence Murphy**
Eric Cox
Catherine Farrell
Tamara Goldfarb
Diana Haddad
John Jackson
Vinita Joardar
Kelly McGarvey
Michael Murphy

Shashi Pujar
Bhanu Rajput
Sanjida Rangwala
Lillian Riddick
Barbara Robbertse
Brian Smith-White
Pooja Strope
David Webb

### RefSeq Developers

Alex Astashyn
Olga Ermolaeva
Vamsi Kodali
Craig Wallin

### Annotation Pipeline

**Francoise Thibaud-Nissen**
Paul Kitts
Mike Dicuccio
Wratko Hlavina
Avi Kimchi

Jinna Choi
Boris Kiryutin
Patrick Masterson
Eyal Mozes
Anton Perkov
Dan Rausch
Robert Smith
Alexandre Souvorov

### GDV/Remap/GBench

Valerie Schneider
Peter Meric
Nathan Bouk
Hsiu-Chuan Chen
Cliff Clausen
Anatoliy Kuznetsov

### A cast of thousands

Ken Katz
Michael Ovetsky
Lukas Wagner
Andrei Shkeda
Donna Maglott
Kim Pruitt
Jim Ostell

Watch NCBI News for updates!
http://www.ncbi.nlm.nih.gov/news/
https://www.youtube.com/user/NCBINLM

NIH) U.S. National Library of Medicine
National Center for Biotechnology Information

# NCBI Genome Resources Workshop

Monday January 13, 2020, 12:50 – 3:00 pm, Pacific Salon 1

| Time | Topic |
|---|---|
| 12:55 – 1:15 | NCBI Wants Your Sequence Data! How Do I Get It There? <br> *Ilene Mizrachi* |
| 1:15 – 1:35 | Annotation of Eukaryote Genomes at NCBI <br> *Jinna Hoffman* |
| 1:35 – 1:55 | Accessing Homologous Gene Datasets at NCBI <br> *Nuala O'Leary* |
| 1:55 – 2:15 | The New PubMed Is Here! <br> *Kathi Canese* |
| 2:15 – 2:35 | Taxonomy Lookup; Data Retrieval: How to Find and Stream Genomic Data in the Cloud! <br> *Ben Busby* |

Visit NCBI Booth **321**      Contact us **info@ncbi.nlm.nih.gov**

NIH〉U.S. National Library of Medicine
National Center for Biotechnology Information

Watch NCBI News for updates!
http://www.ncbi.nlm.nih.gov/news/
https://www.youtube.com/user/NCBINLM

# Also from NCBI!

| Day | Time | Topic |
| --- | --- | --- |
| Saturday | PENDING<br>*Royal Palm Salon 3-4* | PENDING<br>*Aquaculture* |
| Sunday | 10:30 am<br>*Town & Country* | **Stand-Alone PGAP: The NCBI Open-Source Pipeline for the Annotation of Prokaryotic Genomes**<br>*Computational Gene Discovery* |
| Sunday | 12:25 pm<br>*San Diego* | **Genomic Resources for Agricultural Animals at NCBI**<br>*Cattle/Sheep/Goat 2* |
| Sunday | 1:42 pm<br>*Pacific Salon 1* | **NCBI RefSeq Resources for Plant Genomics**<br>*Functional Genomics* |
| Monday | 12:50 pm – 3:00 pm<br>*Pacific Salon 1* | **NCBI Genome Resources Workshop** |
| Tuesday | 11:10 am<br>*California* | **NCBI BLAST: Enhanced Web Usability through New Result Page and Effective Genomic Data Access**<br>*Digital Tools and Resources Session 3* |
| Wednesday | 11:50 am<br>*California* | **Federated Cloud Access to Datasets through Indexing and/or Graphs!**<br>*Digital Tools and Resources Session 4* |